



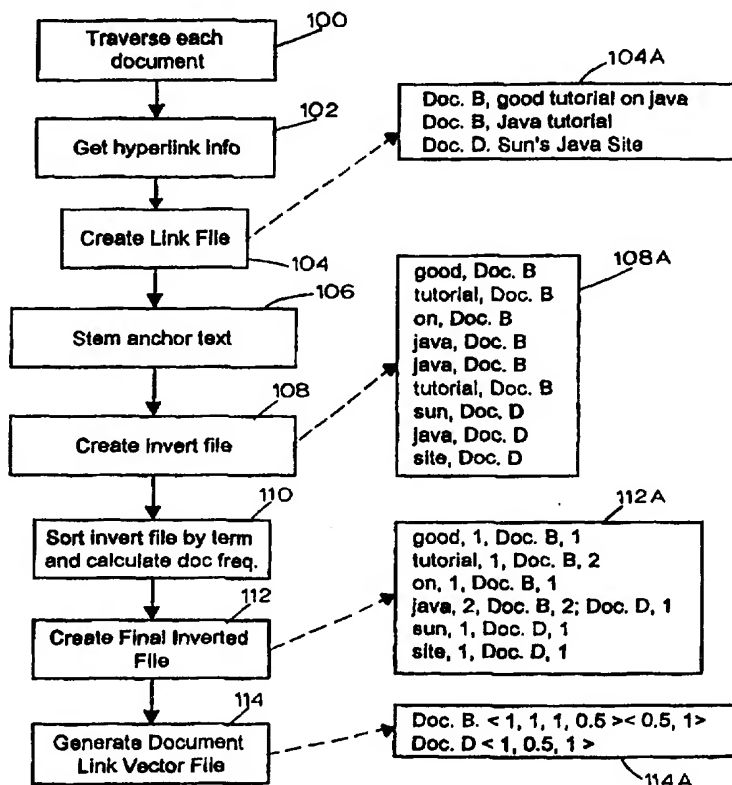
INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification ⁶ : G06F 17/30		A1	(11) International Publication Number: WO 97/49048
			(43) International Publication Date: 24 December 1997 (24.12.97)
(21) International Application Number: PCT/US97/10191 (22) International Filing Date: 17 June 1997 (17.06.97) (30) Priority Data: 08/664,565 17 June 1996 (17.06.96) US 08/794,425 5 February 1997 (05.02.97) US (71) Applicant: IDD ENTERPRISES, L.P. [US/US]; Suite 1810, 2 World Trade Center, New York, NY 10048 (US). (72) Inventor: LI, Yanhong; 191 Spruce Mill Lane, Scotch Plains, NJ 07076 (US). (74) Agent: GERSTEIN, Robert, M.; Marshall, O'Toole, Gerstein, Murray & Borun, 6300 Sears Tower, 233 S. Wacker Drive, Chicago, IL 60606-6402 (US).		(81) Designated States: AL, AM, AT, AU, AZ, BB, BG, BR, BY, CA, CH, CN, CZ, DE, DK, EE, ES, FI, GB, GE, HU, IL, IS, JP, KE, KG, KP, KR, KZ, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, TJ, TM, TR, TT, UA, UG, UZ, VN, ARIPO patent (GH, KE, LS, MW, SD, SZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, ML, MR, NE, SN, TD, TG). Published <i>With international search report.</i> <i>Before the expiration of the time limit for amending the claims and to be republished in the event of the receipt of amendments.</i>	

(54) Title: HYPERTEXT DOCUMENT RETRIEVAL SYSTEM AND METHOD

(57) Abstract

A search engine for retrieving documents pertinent to a query indexes documents in accordance with hyperlinks pointing to those documents. The indexer traverses the hypertext database and finds hypertext information including the address of the document the hyperlinks point to and the anchor text of each hyperlink. The information is stored in an inverted index file, which may also be used to calculate document link vectors for each hyperlink pointing to a particular document. When a query is entered, the search engine finds all document vectors for documents having the query terms in their anchor text. A query vector is also calculated, and the dot product of the query vector and each document link vector is calculated. The dot products relating to a particular document are summed to determine the relevance ranking for each document.



FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AL	Albania	ES	Spain	LS	Lesotho	SI	Slovenia
AM	Armenia	FI	Finland	LT	Lithuania	SK	Slovakia
AT	Austria	FR	France	LU	Luxembourg	SN	Senegal
AU	Australia	GA	Gabon	LV	Latvia	SZ	Swaziland
AZ	Azerbaijan	GB	United Kingdom	MC	Monaco	TD	Chad
BA	Bosnia and Herzegovina	GE	Georgia	MD	Republic of Moldova	TG	Togo
BB	Barbados	GH	Ghana	MG	Madagascar	TJ	Tajikistan
BE	Belgium	GN	Guinea	MK	The former Yugoslav Republic of Macedonia	TM	Turkmenistan
BF	Burkina Faso	GR	Greece	ML	Mali	TR	Turkey
BG	Bulgaria	HU	Hungary	MN	Mongolia	TT	Trinidad and Tobago
BJ	Benin	IE	Ireland	MR	Mauritania	UA	Ukraine
BR	Brazil	IL	Israel	MW	Malawi	UG	Uganda
BY	Belarus	IS	Iceland	MX	Mexico	US	United States of America
CA	Canada	IT	Italy	NE	Niger	UZ	Uzbekistan
CF	Central African Republic	JP	Japan	NL	Netherlands	VN	Viet Nam
CG	Congo	KE	Kenya	NO	Norway	YU	Yugoslavia
CH	Switzerland	KG	Kyrgyzstan	NZ	New Zealand	ZW	Zimbabwe
CI	Côte d'Ivoire	KP	Democratic People's Republic of Korea	PL	Poland		
CM	Cameroon	KR	Republic of Korea	PT	Portugal		
CN	China	KZ	Kazakhstan	RO	Romania		
CU	Cuba	LC	Saint Lucia	RU	Russian Federation		
CZ	Czech Republic	LI	Liechtenstein	SD	Sudan		
DE	Germany	LK	Sri Lanka	SE	Sweden		
DK	Denmark	LR	Liberia	SG	Singapore		
EE	Estonia						

- 1 -

HYPERTEXT DOCUMENT RETRIEVAL SYSTEM AND METHOD**FIELD OF INVENTION**

The present invention relates to hypertext document retrieval,
and more particularly to systems and methods of searching databases
5 distributed over wide-area networks such as the World Wide Web.

BACKGROUND OF THE ART

A hypertext is a database system which provides a unique
and non-sequential method of accessing information using nodes and links.
Nodes, i.e. documents or files, contain text, graphics, audio, video,
10 animation, images, etc. while links connect the nodes or documents to
other nodes or documents. The most popular hypertext or hypermedia
system is the World Wide Web, which links various nodes or documents
together using hyperlinks, thereby allowing the non-linear organization of
text on the web.

15 A hyperlink is a relationship between two anchors, called the
head and the tail of the hyperlink. The head anchor is the destination node
or document and the tail anchor is the document or node from which the
link begins. On the web, hyperlinks are generally identified by
underscoring or highlighting certain text or graphics in a tail anchor

- 2 -

document. When a user reviewing the tail document "clicks on" the highlighted or "anchor-text" material, the hyperlink automatically connects the user's computer with or "points to" the head anchor document for that particular hyperlink.

5 A hypertext system generally works well when a user has already found a tail document pertaining to the subject matter of interest to that user. The hyperlinks in the tail document are created by the author of the document who generally will have reviewed the material in the head documents of the hyperlinks. Thus, a user clicking on a hyperlink has a
10 high degree of certainty that the material in the head document has some pertinence to the anchor text in the tail document of the hyperlink.

 As the popularity of the Internet and the Web has grown, the ability to find relevant documents has become increasingly difficult. If a user is unable to find a first document pertaining to the subject matter of
15 interest, the user will of course not be able to use hyperlinks to find additional pertinent documents. Moreover, the location of a single relevant document may not lead to other documents if the author of the relevant document has not created hyperlinks to other relevant web sites. The proliferation of information has, therefore, lead to the development of
20 various search engines which assist users in finding information. Numerous search engines such as Excite, Infoseek, and Yahoo! are now available to users of the Web.

- 3 -

Search engines usually take a user query as input and attempt to find documents related to that query. Queries are usually in the form of several words which describe the subject matter of interest to the user. Most search engines operate by comparing the query to an index of a

5 document collection in order to determine if the content of one or more of those documents matches the query. Since most casual users of search engines do not want to type in long, specific queries and tend to search on popular topics, there may be thousands of documents that are at least tangentially related to the query. When a search engine has indexed a large

10 document collection, such as the Web, it is particularly likely that a very large number of documents will be found that have some relevance to the query. Most search engines, therefore, output a list of documents to the user where the documents are ranked by their degree of pertinence to the query and/or where documents having a relatively low pertinence are not

15 identified to the user. Thus, the way in which a search engine determines the relevance ranking is extremely important in order to limit the number of documents a user must review to satisfy that user's information needs.

Almost all ranking techniques of search engines depend on the frequency of query terms in a given document. When other related

20 factors are the same, the higher a term's frequency in a given document, the higher the relevance score of this document to a query including that term. Factors other than term frequency, such as such document frequency, i.e. how many documents contain the term, may also be taken

- 4 -

into account in determining a relevance score. Once the various factors such as term frequency or document frequency have been determined for a particular query, various models such as the vector space model, probabilistic model, fuzzy logic models, etc. are used to develop a numerical relevance ranking. See, Harman, D., "Ranking Algorithms," Chapter 14, *Information Retrieval*, (Prentice Hall, 1992).

For instance, in the vector space model, a user query Q is represented as a vector where each query term (qt) is represented as a dimension of a query vector.

10
$$Q = \langle qt_1, qt_2, \dots, qt_m \rangle$$

Documents in the database are also represented by vectors with each term or key word (dt) in the document represented as a dimension in the vector.

$$D = \langle dt_1, dt_2, \dots, dt_n \rangle$$

The relevance score is then calculated as the dot product of Q and D.

15 The calculation of the value of each dimension for vectors Q or D may be weighted in a variety of ways. The most popular term-weighting formula is:

$$\text{Weight}(t) = \text{TF} * \text{IDF}_t$$

where TF is the term frequency of a given term in a document or query, and IDF_t is the inverse document frequency of the term. The inverse document frequency is the inversion of how many documents in the whole document collection contain the term, i.e.:

20

- 5 -

$$IDF_t = \frac{1}{DF_t}$$

Using an inverse document frequency insures that junk words such as "the," "of," "as," etc. do not have a high weight. In addition, when a query uses multiple terms, and one of those terms appears in many documents, using an IDF weighting gives a lower ranking to documents
5 containing that term, and a higher ranking to document containing other terms in the query.

There are normalized versions of term weighting, which take into account the length of a document including a particular term. The assumption made is that the more frequently a term appears in a document
10 for a given amount of text, the more likely that document is relevant to a query including that term. That assumption may not be true, however, in many cases. For example, if the query is "Java tutorial," a document (call it J), which contains 100 lines with each line consisting of just the phrase "Java tutorial," would get a very high relevance score and would be output
15 by a search engine as one of the most relevant documents to the user. That document, however, would be useless to the user since it provides no information about a "Java tutorial." What the user really needs is a good tutorial for the Java programming language such as found on Sun's Java tutorial site (<http://Java.sun.com/tutorial>). Unfortunately, the phrase "Java
20 tutorial" does not occur 100 times on Sun's site, and therefore most search

- 6 -

engines would incorrectly find Sun's site to be less pertinent, and thus have a lower relevance ranking, than Document J.

Documents such as Document J might not be included in a traditional database because each document in a traditional database is
5 selected or authored for its content rather than the repetition of certain key words. On the Web, where anyone can be a publisher, there is no one to select or screen out document such as J. In fact, some people intentionally draft their documents so that the documents will be retrieved on the top of a ranked list output by search engines that take into account term frequency
10 or normalized term frequency. For instance, a Web site may be designed so that the text for the first five lines includes the word "sex." The Web site may be of low quality or have nothing to do with sex, but a search engine can be fooled into ranking the site highly because of the high frequency of the word "sex" in the site.

15 Length normalization may also have other problems in a hypertext environment. Documents containing media other than text may make it difficult to accurately calculate the relevant length of a document.

Traditional search engines using key words also may not retrieve relevant documents containing synonyms of those key words.
20 Thus, many search engines may need an extensive thesaurus, which may be too expensive or difficult to build, in order to find a document containing the word "attorney" when the user includes only the word "lawyer" in a query. Traditional search engines also cannot find relevant documents

- 7 -

which are in a language other than the language of the query entered by the search engine user. Translation tools are a possible solution, but they may be difficult and expensive to build.

In addition, traditional search engines may be unable to
5 identify non-textual material which is relevant to a query. For instance, a Web site containing pictures of Mozart or examples of Mozart's music may not be deemed relevant by a search engine when that search engine can only search for the word "Mozart" within the text of documents.

SUMMARY OF THE INVENTION

10 A method of indexing documents includes obtaining a list of hyperlinks pointing to each document, where each hyperlink includes one or more terms. Each document is indexed with the terms in the hyperlink pointing to that document. A number of hyperlinks, each containing a particular term, may point to a document. The number of hyperlinks
15 containing that particular term pointing to the document is indexed with that document.

A particular term may appear in hyperlinks pointing to a number of documents, and the number of documents having the particular term in hyperlinks pointing to those documents is indexed with that term.
20 Indexing may include creating a file listing each term, the number of documents having that term in hyperlinks pointing to those documents, a

- 8 -

document identifier for each document having that term in hyperlinks pointing that document, and the number of hyperlinks containing that term pointing to each identified document.

The number of documents having a particular term in
5 hyperlinks pointing to those documents may be indexed with a document identifier for each document having the particular term in a hyperlink pointing to that document. The indexing of a particular term in a hyperlink pointing to a document may be with the inverse of the number of documents having the particular term in hyperlinks pointing to those
10 documents.

A term may appear a number of times in a hyperlink pointing to a document, and the number of times each term appears in a hyperlink is indexed with the document pointed to by the hyperlink.

The terms may be stemmed words. The method of the
15 present invention may be performed on an apparatus and may be stored as a computer-readable set of instructions.

In accordance with another aspect of the present invention, a method of ranking documents is based on the document's relevance to a query where the query has at least one term, and where hyperlinks contain
20 terms and point to corresponding documents. The method includes comparing the words in the query to the words in a hyperlink to obtain a relevance ranking for each hyperlink, and summing the relevance rankings

- 9 -

for each hyperlink pointing to a particular document to obtain a summed relevance score for that document.

The query may be represented by a query vector where the query vector contains a dimension for each term in the query. Each document may be represented by document link vectors for each hyperlink pointing to the document, where each document link vector contains a dimension for each term in the corresponding hyperlink pointing to that document. Comparing the words in the query to the words in the hyperlinks includes calculating the dot product of the query vector with the document link vector for that hyperlink. Summing the relevance ranking for each hyperlink pointing to a document includes summing the dot products obtained using the document link vectors for a particular document to obtain the summed relevance score for that document. The summed relevance scores may then be compared to obtain a ranking of documents.

The dimension for a term in a query vector may be related to the inverse of the number of documents having a respective hyperlink containing that term pointing to those documents. Similarly, the dimension for a term in a document link vector may be related to the inverse of a number of documents having a respective hyperlink containing that term pointing to those documents.

Other features and advantages are inherent in the hypertext document retrieval system and method claimed and disclosed or will

- 10 -

become apparent to those skilled in the art from the following detailed description in conjunction with the accompanying drawings.

BRIEF DESCRIPTION OF THE DRAWINGS

Fig. 1 is a block diagram of a distributed computer network
5 including a hypertext retrieval system of the present invention;

Fig. 2 is a block diagram of an indexing and retrieval system
of the present invention;

Fig. 3 is a diagram of two hypertext documents;

Fig. 4 is an example of a hypertext document system
10 including representation of hyperlinks between those documents;

Fig. 5 is a flow chart of an indexing process of the present
invention; and

Fig. 6 is a flow chart of a retrieval process of the present
invention.

15 DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENT

Fig. 1 is a block diagram of a typical distributed hypertext system including a client computer 10 connected to server computers 12, 13, 14, 15, and 16. Although the client computer 10 is shown connected directly to server 12, it may be connected indirectly to server 12 through a

- 11 -

service provider or through any one or more of the other servers. Servers 13, 14, 15, and 16 include files of documents 17, 18, 19, and 20, respectively. Files 17, 18, 19, and 20 contain documents available to users of the network. Server 12 includes an index file 21 as discussed in more
5 detail below. The server computer 12 traverses the network looking for all hypertext documents residing in the files 17-20 of the other server computers 13-16 in order to build the index file 21.

Fig. 2 describes the general structure of an indexing and retrieval system 30 of the present invention. A user from outside the
10 system 30 inputs a query 32 through a user interface 34, which will typically reside on the user's computer, such as a client computer 10 (Fig. 1). The user's query is then transmitted through the network to the indexing and retrieval system 30, which generally resides on a server, such as server 12 (Fig. 1). The system 30 includes a retrieval engine 36, index
15 files 38, and an index engine 40. The operation of the retrieval engine 36 and index engine 40 and the creation of the index files 38 are described below. The index engine 40 creates the index files 38 by traversing a document database 42, such as that found on the World Wide Web. The document database 42 might include files 17-20 (Fig. 1). The index files
20 38 created by the index engine 40 may take various forms in accordance with the present invention, but may include a link file 44, an inverted file 46, and a document vector file 48, all of which are described in detail below. The retrieval engine 36 uses the index files 38 in order to

- 12 -

determine a relevance ranking for the documents, and outputs search results at 49 through the user interface 34.

Fig. 3 is a diagram of Document A and Document B, which are the tail anchor and head anchor, respectively, of the hyperlink

5 represented by the arrow 50. The Document A has an address "URL1" and Document B has an address "URL2." The addresses may be in the form of a uniform resource locator (URL), which is a type of uniform resource identifier (URI) for head and tail anchor addresses. URL's are typically in the format such as:

10 `http://www.w3.org/hypertext/book.html`

Optionally, the URL may be followed by the pound symbol and a sequence of characters called a fragment identifier in order to identify a fragment within a document, i.e.:

`http://www.w3.org/hypertext/book.html#Chapter1`

15 Document A has a title 52, an abstract 54, and text or media 56.

Similarly, Document B has a title 58, an abstract 60, and text or media 62.

The text or media may contain anchor text such as anchor text 64 in Document A. Document A also contains a command 66, which serves as the instructions for the hyperlink 50. The representation of
20 command 66 of the hyperlink 50 is shown in hypertext markup language (HTML) and includes the command "href" and then identifies the address of the head anchor, in this case, the address of Document B "URL2." The command 66 then includes the statement "good tutorial on Java," which

- 13 -

identifies the anchor text of the hyperlink 50. By identifying the phrase "good tutorial on Java" as the anchor text in the command 66, that phrase is thereby underlined in the text 56 of Document A. When text such as anchor text 64 is underlined, it alerts a reader of Document A to the
5 existence of the hyperlink. When a user then clicks on the anchor text 64, the command 66 points to Document B, thereby instructing the user's computer to send a message to the address URL2, requesting a copy of Document B.

The author of Document A must, of course, create the
10 command 66 and identify the anchor text 64. Generally, authors of such documents will describe, in that author's opinion, the head anchor document (in this case Document B) with the words of the anchor text (in this case, anchor text 64). Therefore, if there are many authors like the author of Document A that make link commands to document B using the
15 anchor text 64, then a user looking for a Java tutorial is highly likely to be interested in the information in Document B.

Fig. 4 is a representation of a simple hypertext system having only four documents, Documents A, B, C, and D. The system shown in Fig. 4 has only three hyperlinks, hyperlink 50, also shown in
20 Fig. 3, and hyperlinks 68 and 70. The anchor text "good tutorial on Java" in Document A is the tail for the hyperlink from Document A to Document B, as shown in Fig. 3. Document C contains two sets of anchor text "Java tutorial" and "Sun's Java site." The anchor text 72 in Document C points

- 14 -

to Document B through the hyperlink 68. The anchor text 74 points to Document D through the hyperlink 70. The hypertext system shown in Fig. 4 will be used below in describing the hypertext system including the index engine, the retrieval engine, and the index files created by the index engine.

Fig. 5 describes the operation of the index engine 40 of Fig. 2. At block 100, the index engine traverses each document in the database. Traversing the database can be accomplished in a variety of ways, but usually using a so-called "spider" program. See, Cheong, F.C. *Internet Agents: Spiders, Wanderers, Brokers, and Bots*, (McMillan, 1997). Spider programs begin by obtaining various URL addresses and send messages to those addresses requesting the documents located at the addresses. Those addresses may identify a server, a document stored in files on that server, or groups of documents. Upon obtaining the document or documents identified by the URL, a spider program then reviews those documents looking for hyperlink commands identifying additional addresses. The spider program records those addresses and then seeks the documents residing at those addresses.

While traversing each document in block 100, the system also obtains hyperlink information at block 102 regarding each document. Such hyperlink information might include the URL of the document, the words in the anchor text of the hyperlink in the document, and the URL of any document pointed to by a hyperlink having that anchor text. The

- 15 -

system may also collect a variety of information about the document including its title and possibly the text of the document. The system may also create an abstract, if desired.

At block 104, the system creates one or more link files
5 where entries in the files have a format:

< doc.ID, anchor-text >

where doc.ID is an identifier for each head document of a hyperlink having the corresponding anchor text. The doc.ID may be in the form of a URL or may be another identifier which is indexed in some manner with the
10 document's URL. Box 104A is an example of a link file, as referred to in Fig. 2, created for the database of the documents shown in Fig. 4. Since the database in Fig. 4 has three hyperlinks, there are three entries in file 104A. The system may also store the number of times a term appears in anchor text for a particular link. In the examples shown, each term only
15 appears once in a particular link.

Although Fig. 5 shows that traversing of documents in block 100 occurs before link files are created at block 104, it is possible for some link files to be created prior finishing traversing all documents in the database. In fact, once the database has been entirely traversed, it may be
20 desirable to update the link files and other index files by retraversing documents in order to determine if any additional documents have been added to the database, or if any hyperlinks have been added to the documents.

- 16 -

At block 106, the anchor text for the various hyperlinks may be stemmed. Stemming is a process of reducing the words from various morphological forms to a simplified stem. During stemming, words are usually made case-insensitive, e.g. "Tutorial" and "tutorial" are considered
5 the same. "Sun's" will stem to "Sun," "documents" will stem to "document," etc.

Control then passes to block 108, which creates an invert file with entries in the format of:

< term, doc. >

10 where term is a word extracted from the anchor text of a hyperlink and doc. is the identifier for the head document of that hyperlink. An invert file as created in block 108 is shown in file 108A. Since the anchor text "good tutorial on Java" has four words, that hyperlink results in four entries in file 108A.

15 At block 110, the invert file is sorted by term, and the document frequency is calculated. The document frequency is defined as the number of documents having a particular term in anchor text of hyperlinks pointing to those documents. For instance, in the database of Fig. 4, the term "Java" appears in the anchor text of three hyperlinks,
20 where those three hyperlinks point to a total of two different documents. Therefore, the document frequency for the term "Java" is two. The term "good" appears in only one hyperlink that points to only one document, so the document frequency for the term "good" is one.

- 17 -

Control next passes to block 112, which creates final invert file as shown in 112A. Entries in the final invert file are in the format:

$\langle \text{term}, \text{DF}, \text{doc1}, \text{lf1}, \text{doc2}, \text{lf2}, \dots, \text{doci}, \text{LFi} \rangle$

where "term" is a term in the anchor text, DF is the document frequency
5 for that term, doci is the document identifier for Document i, and LFi is the link term frequency for doci. Link term frequency is defined as the number of hyperlinks pointing to doci whose anchor text consists of the particular term. For example, the term "good" appears in only one
10 hyperlink that points to Document B, so the link term frequency of the term "good" for Document B is one. The term "Java" appears in two hyperlinks that point to Document B, so the link term frequency of "Java" for Document B is two. One embodiment of the retrieval engine of the present invention will depend on this file to find documents related to a user query.

15 The index engine at box 114 may also generate a document link vector file where entries in the document link vector file are in the format of:

$\text{doc.id}, v_1, v_2, \dots, v_i$

where doc.id is the identifier for a particular document, and v_i is a vector
20 representation of a hyperlink found in the link file. Each vector v_i will be in the format of:

$\langle w(t_1), w(t_2), \dots, w(t_i) \rangle$

- 18 -

where $w(t_i)$ is the weight of term i in a given anchor text for the hyperlink represented by the vector. The dimension of each document link vector ($w(t_i)$) is calculated by $TF_i * IDF$, where TF_i is the term frequency of term i , i.e. how many times a term appears in the given anchor text, and IDF is the invert document frequency ($1/DF$) for the term to which the particular dimension in the link vector pertains. It may be desirable to divide the document frequency by the total number of documents to obtain a normalized document frequency when calculating the dimensions. It may also be desirable to use the logarithm of the inverse document frequency when calculating dimensions.

File 114A is an example of a document link vector file which has been generated at block 114. Since Document B has two hyperlinks pointing to it, there are two vectors for Document B entered in file 114, along with the identifier of Document B. Since the anchor text of the first hyperlink pointing to Document B has four distinct words "good tutorial on Java," the first vector for Document B has four dimensions. Since the second hyperlink pointing to Document B has only two words in the anchor (Java, tutorial), the second vector indexed with Document B has only two dimensions.

As described below, the document link vector file 114A is used in calculating the relevance score with respect to a particular query. Instead of creating document link vector files automatically, it may be desirable to create document link vector files only upon receipt of a query.

- 19 -

Thus, the only entries in the link vector files which need to be created are those pertaining to documents having query terms in the anchor text of hyperlinks pointing to those documents.

In the first vector for Document B, the first three dimensions
5 are "one" since the terms "good," "tutorial," and "on" only appear in anchor text pointing to one document, and they only appear once in the anchor text. Thus:

$$TF*IDF = 1*1=1.$$

The term "Java," however, has a term frequency of one and document
10 frequency of two, and therefore has an inverse document frequency of .5. Thus, $TF*IDF$ for "Java" is .5, making the last dimension in the first vector for Document B equal to .5. The remaining dimensions in the second vector for Document B and the vector for Document D are also calculated according to the $TF*IDF$ formula.

15 The link file 104A, the invert file 108A, the final invert file 112A, and the document link vector file 114 are all considered index files as shown in Fig. 2. Although the files as shown in Fig. 5 are preferred, there are many indexing techniques which can be used with a system of the present invention, which rely on anchor text and link frequency in order to
20 index documents. For instance, the files may be compressed or have a variety of relational structures for the data within files or between files.

Referring now to Fig. 6, the retrieval process achieves relevance ranking by using the vector space model and link vector voting.

- 20 -

The process begins at box 120 with the input of a user query as shown in file 120A. At box 122, the system then searches the inverted file or final inverted file and, at box 124, finds all documents indexed with the query terms. A document may be related to the query if that document has a
5 hyperlink pointing to it, where the hyperlink includes a query term in its anchor text. As shown in box 124A, the system has located two documents, Document B and Document D, each of which has one or more of the terms in the query in anchor text of hyperlinks pointing to those documents.

10 Control next passes to box 126 where the system finds document link vectors for each document identified in box 124A. The document link vectors are contrasted with conventional document vectors which are based on the content of each document. The system may find the document link vectors by simply going to the document link vector file
15 114 (Fig. 5) or may create the document link vectors from the invert file and link file. Box 126A shows the document link vectors, along with the anchor text, for each hyperlink pointing to a document related to the query.

 While obtaining the document link vectors, the system, at box 128, also creates a query vector as shown in box 128A. The
20 dimensions in the query vector are equal to $TF_q * IDF$ for each term in the query, where TF_q is the term frequency or number of times the term appears in the query. IDF is the inverse document frequency for a term as calculated in box 110 of Fig. 5. The TF_q is one for both "Java" and

- 21 -

"tutorial" in the query. The IDF as previously calculated in box 110 of Fig. 5 for "Java" is .5 and as calculated for "tutorial" is one.

Once the query vector and all relevant document link vectors have been found or calculated, control passes to block 130 to calculate the
5 relevance scores for each document. The relevance score is calculated by finding the dot product of each document link vector with the query vector.

A dot product for vectors $\langle a, b, c \rangle$ and $\langle d, e, f \rangle$ is defined as:

$$\frac{a*d + b*e + c*f}{\sqrt{a^2+b^2+c^2}\sqrt{d^2+e^2+f^2}}$$

10 If two vectors do not have the same dimensions, a zero is entered for each dimension which is not present in that vector. For instance, the first vector for Document B is represented as:

$$\langle 1, 1, 1, 0.5 \rangle.$$

In such an instance, the query vector would be represented as:

15 $\langle 0, 1, 0, .5 \rangle$

so that the dimensions representing "tutorial" in each vector and "Java" in each vector match up. The dot product of the query vector with the first document link vector for Document B would then be calculated as follows:

20
$$\frac{0 \times 1 + 1 \times 1 + 0 \times 1 + .5 \times .5}{\sqrt{1^2+1^2+1^2+.5^2}\sqrt{1^2+.5^2}} = .620$$

A similar calculation for the second vector for Document B would lead to a dot product of 1.

- 22 -

At box 131, the dot products for all document link vectors pertaining to a particular document are summed to obtain a "vote" or summed score for a particular document. The summed relevance score for Document B is the sum of the dot products for each document link vector relating Document B, which equals 1.620. A similar calculation can be made by finding the dot product of the query vector with the only document link vector for Document D, which equals 0.149.

At box 132, sorted results are output as shown in box 132A. The results are sorted so that the documents having higher summed relevance rankings are listed above those with lower rankings. Instead of listing all documents having a non-zero relevance score, it may be desirable to only list a pre-set number, i.e. the top 100 documents, or to only list those documents having a relevance score above a certain threshold.

The process described herein can be performed on a number of apparatus, including a Sun Sparc Station with a Solaris operating system. The process may be stored in memory on the computer system as a set of instructions. The set of instructions may also be stored on a computer-readable memory such as a disk, and the instructions can be transmitted from one computer to another over a network.

In the example described, no hyperlinks point to Document A or C, so each of their relevance scores is zero, even though both Document A and Document C contain the words in the query, "Java" and "tutorial." A conventional index and retrieval engine could be used in

- 23 -

combination with the hyperlinked based index and retrieval system of the present invention. This combination might be used in the case of a link-based relevance score tie, or merely to supplement the link-based information. For instance, suppose the relevance scores for Document A and C are 0.6 and 0.8, respectively, based on conventional and relevance
5 ranking. The final relevance ranking for the query utilizing the conventional ranking to break the tie of the link-based ranking would be Document B, Document D, Document C, and Document A.

Another reason to use combination ranking may be when
10 there are too few hyperlinks (such as only one link) pointing to a document. In such a case, the relevance score based upon the one link may not be accurate, so a threshold can be set for the link-based relevance score. If the link-based relevance score is lower than the threshold, other means of relevance ranking may be used or combined with the link-based
15 relevance score.

Because the index files of the present invention use only hyperlink information, relevance ranking does not depend on the words appearing in documents themselves, or, if used in combination with conventional relevance ranking do not depend solely on words appearing in
20 the documents. Instead, the relevance ranking depends on descriptions of those documents in the anchor text of hyperlinks pointing to the documents. Documents such as Document J described above will not have a high

- 24 -

summed relevance score because authors creating hypertext documents will not include hyperlinks in their documents pointing to Document J.

The size of a document is no longer a factor in the relevance ranking, and therefore problems associated with document size can be
5 avoided.

The use of thesauruses may be less important because even if the word "lawyer" never appears in a document titled "California Immigration Attorneys," someone may have created a hyperlink pointing to that document where the anchor text includes the word "lawyer."

10 Images, graphics, and sounds, which are not searchable by conventional information retrieval methods, are searchable if there are hyperlinks pointing to them. Anchor text may also be in the form of images, graphics, etc. so the index engine may substitute other information such as the tail document's title for the non-textual anchor text.

15 Documents in a foreign language may also be retrieved if indexing is performed in accordance with the present invention. If documents written in English contain anchor pointing to the foreign-language documents, the foreign-language documents will receive a relevance score in accordance with the present invention.

20 Thus, when a document database is large enough, as in the case of the World Wide Web, search results are based on a kind of voting, where the description of the content of a document is determined by how others describe the document rather than simply by how the document

- 25 -

describes itself. Thus, in the examples shown above, Sun's Java tutorial site will receive a high summed relevance rank even though the term "Java tutorial" appears only once in the document.

The ranking method based on hyperlinks pointing to a given
5 document can be used to select the most popular documents in a specific field using the feature words or description of that field as the query to the system. By analyzing the link file described in the preferred embodiment, and comparing the different descriptions of hyperlinks pointing to the same document, a system can automatically construct a thesaurus or synonym
10 tool.

The foregoing detailed description has been given for clearness of understanding only, and no unnecessary limitations should be understood therefrom, as modifications would be obvious to those skilled in the art.

- 26 -

Claims

1. A method of indexing documents, the method
2 comprising:
obtaining a list of hyperlinks pointing to each document,
4 wherein each hyperlink includes one or more terms;
indexing each document with the terms in the hyperlinks
6 pointing to that document, wherein a number of hyperlinks, each containing
a particular term, may point to a document; and
8 indexing the number of hyperlinks containing the particular
term pointing to the document with that document.
2. The method of claim 1 wherein:
2 a particular term may appear in hyperlinks pointing to a
number of documents; and
4 the number of documents having the particular term in
hyperlinks pointing to those documents is indexed with that term.

- 27 -

3. The method of claim 2 wherein the indexing
2 comprises creating a file listing:
each term;
4 the number of documents having that term in hyperlinks
pointing to those documents;
6 a document identifier for each document having that term in
hyperlinks pointing to that document; and
the number of hyperlinks containing that term pointing to
each identified document.

4. The method of claim 1 wherein:
2 a particular term may appear in hyperlinks pointing to a
number of documents; and
4 the number of documents having the particular term in
hyperlinks pointing to those documents is indexed with a document
6 identifier for each document having the particular term in a hyperlink
pointing to that document.

5. The method of claim 4 wherein each document having
2 a particular term in a hyperlink pointing to that document is indexed with
an inverse of the number of documents having the particular term in
4 hyperlinks pointing to those documents.

- 28 -

6. The method of claim 1 wherein:
- 2 a term may appear a number of times in a hyperlink pointing
to a document; and
- 4 the number of times each term appears in a hyperlink is
indexed with the document pointed to by the hyperlink.
7. The method of claim 1 wherein the terms are
- 2 stemmed words.
8. An apparatus comprising means for performing the
- 2 method of claim 1.
9. A computer-readable memory device comprising a set
- 2 of instructions for performing the method of claim 1.

- 29 -

10. A method of ranking documents based on the
2 document's relevance to a query, wherein the query comprises at least one
term, and wherein hyperlinks contain terms and point to corresponding
4 documents, the method comprising:

comparing the words in the query to the words in a
6 hyperlink to obtain a relevance ranking for each hyperlink; and
summing the relevance rankings for each hyperlink pointing
8 to a particular document to obtain a summed relevance score for that
document.

11. The method of claim 10 wherein:
2 a number of hyperlinks, each containing a particular term,
may point to a document; and
4 the number of hyperlinks containing the particular term
pointing to the document is indexed with that document.

12. The method of claim 11 wherein:
2 a particular term may appear in hyperlinks pointing to a
number of documents; and
4 the number of documents having a particular term in
hyperlinks pointing to those documents is indexed with that term.

- 30 -

13. The method of claim 12 comprising the creation of a
2 list wherein the list indexes:
each term;
4 the number of documents having hyperlinks pointing to those
documents;
6 a document identifier for each document; and
the number of hyperlinks containing that term pointing to
8 each document.

14. The method of claim 10 wherein:
2 a particular term may appear in hyperlinks pointing to a
number of documents; and
4 the number of documents having the particular term in
hyperlinks pointing to those documents is indexed with a document
6 identifier for each document having the particular term in a hyperlink
pointing to that document.

15. The method of claim 14 wherein each document
2 having a particular term in a hyperlink pointing to that document is indexed
with an inverse of the number of documents having the particular term in
4 hyperlinks pointing to those documents.

- 31 -

16. The method of claim 10 wherein:
2 a term may appear a number of times in a hyperlink pointing
to a document; and
4 the number of times each term appears in a hyperlink is
indexed with the document pointed to by the hyperlink.

17. The method of claim 10 wherein the terms are
2 stemmed words.

18. The method of claim 10 wherein:
2 the query is represented by a query vector wherein the query
vector contains a dimension for each term in the query; and
4 each document is represented by document link vectors for
each hyperlink pointing to the document, wherein each document link
6 vector contains a dimension for each term in the corresponding hyperlink
pointing to that document.

19. The method of claim 18 wherein comparing the words
2 in the query to the words in the hyperlink comprises calculating the dot
product of the query vector with the document link vector for that
4 hyperlink.

- 32 -

20. The method of claim 19 wherein summing the
2 relevance ranking for each hyperlink pointing to a document comprises
summing the dot products obtained using the document link vectors for a
4 particular document to obtain the summed relevance score for that
document.

21. The method of claim 20 wherein the summed
2 relevance scores for each document are compared to obtain a ranking of
documents.

22. The method of claim 18 wherein the dimension for a
2 term in a query vector is related to the inverse of the number of documents
having a respective hyperlink containing that term pointing to those
4 documents.

23. The method of claim 18 wherein the dimension for a
2 term in a document link vector is related to the inverse of the number of
documents having a respective hyperlink containing that term pointing to
4 those documents.

24. An apparatus comprising means for performing the
2 method of claim 10.

- 33 -

25. A computer-readable memory device comprising a set
2 of instructions for performing the method of claim 1.

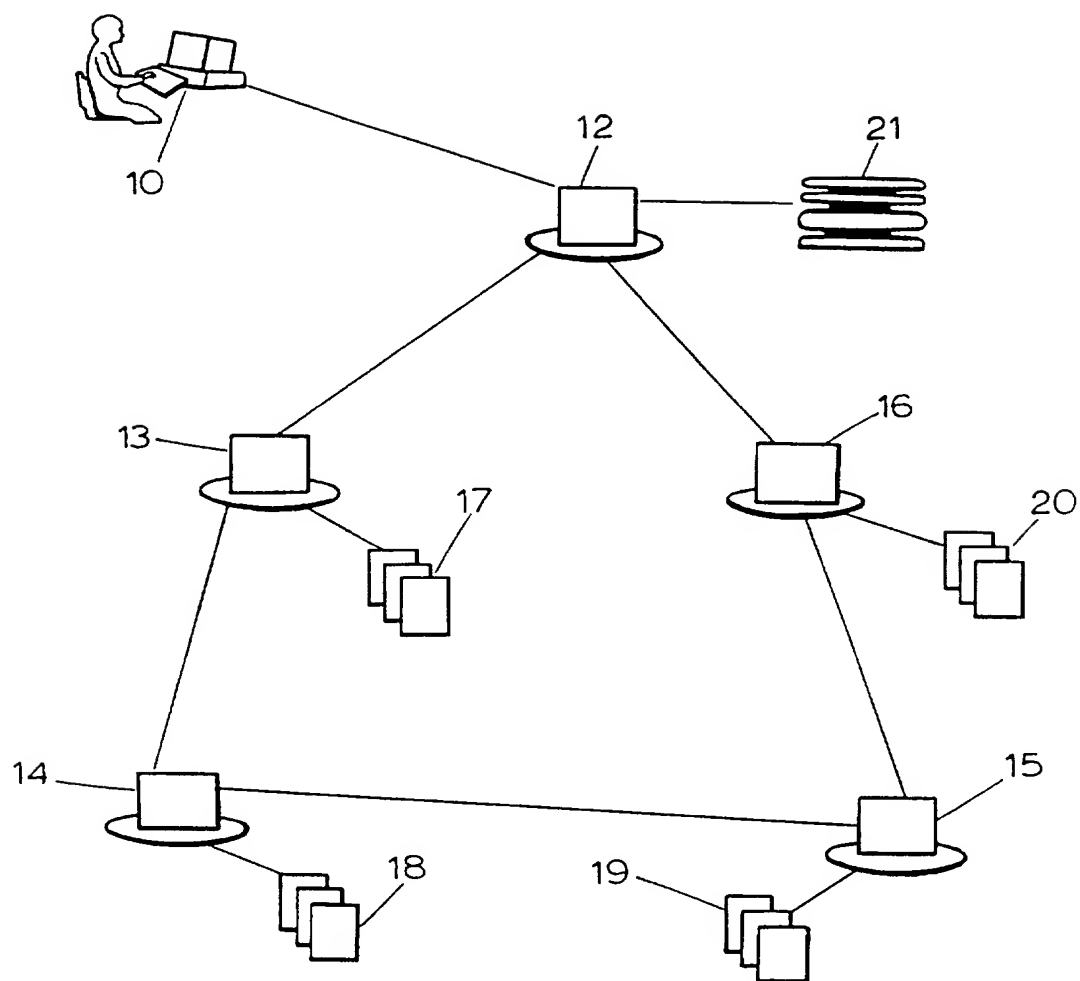


FIG. 1

2/6

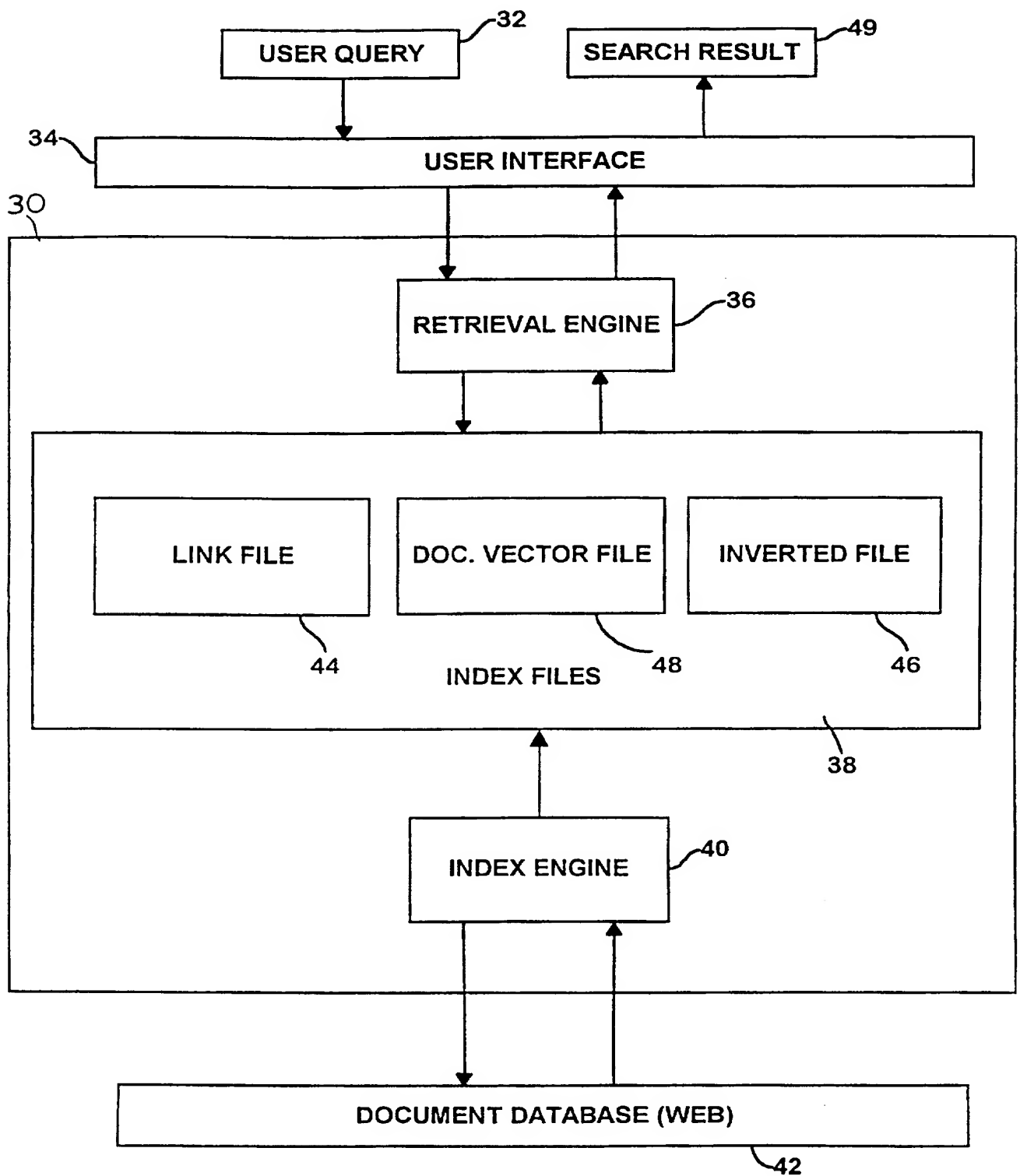


FIG. 2

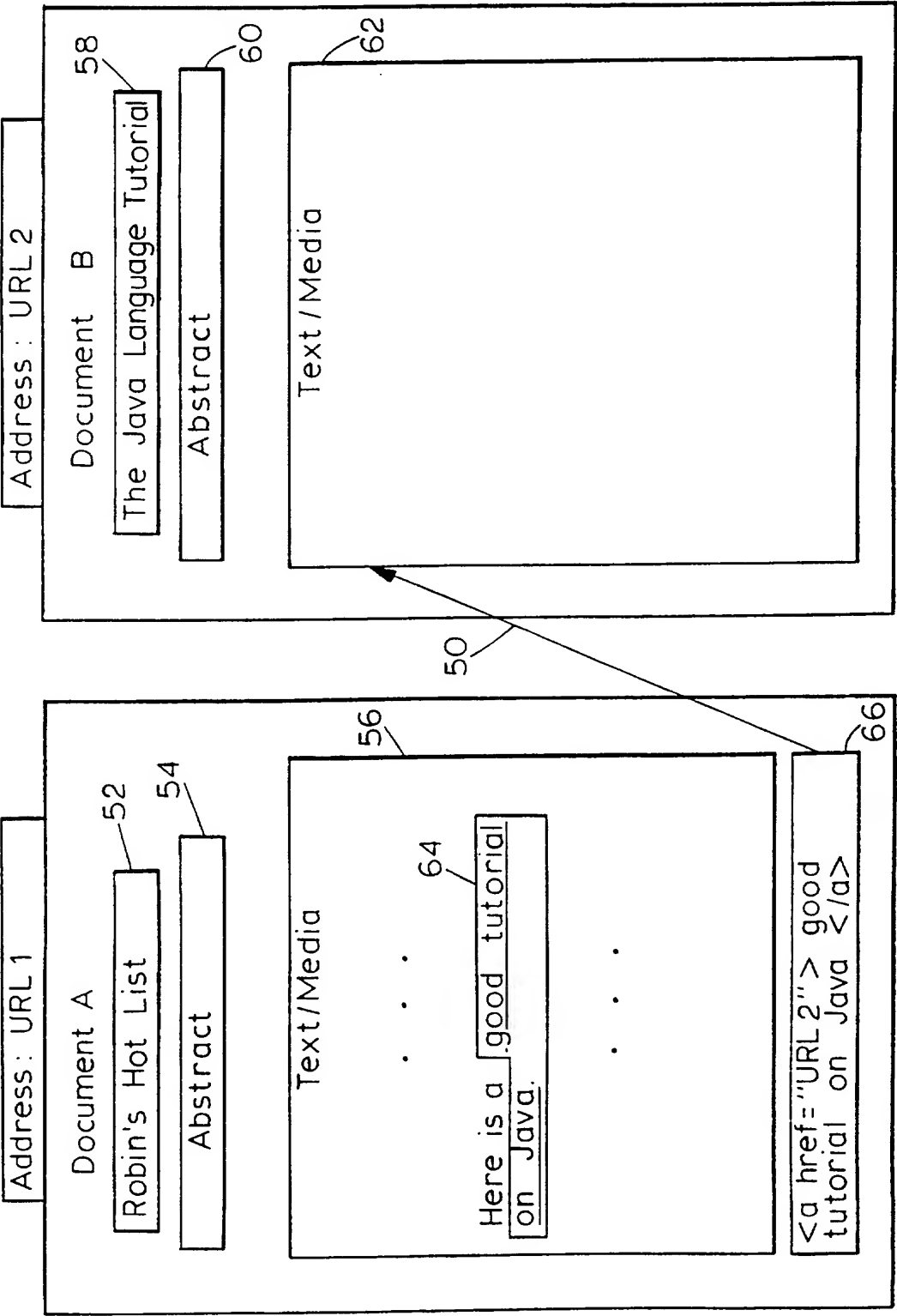


FIG. 3

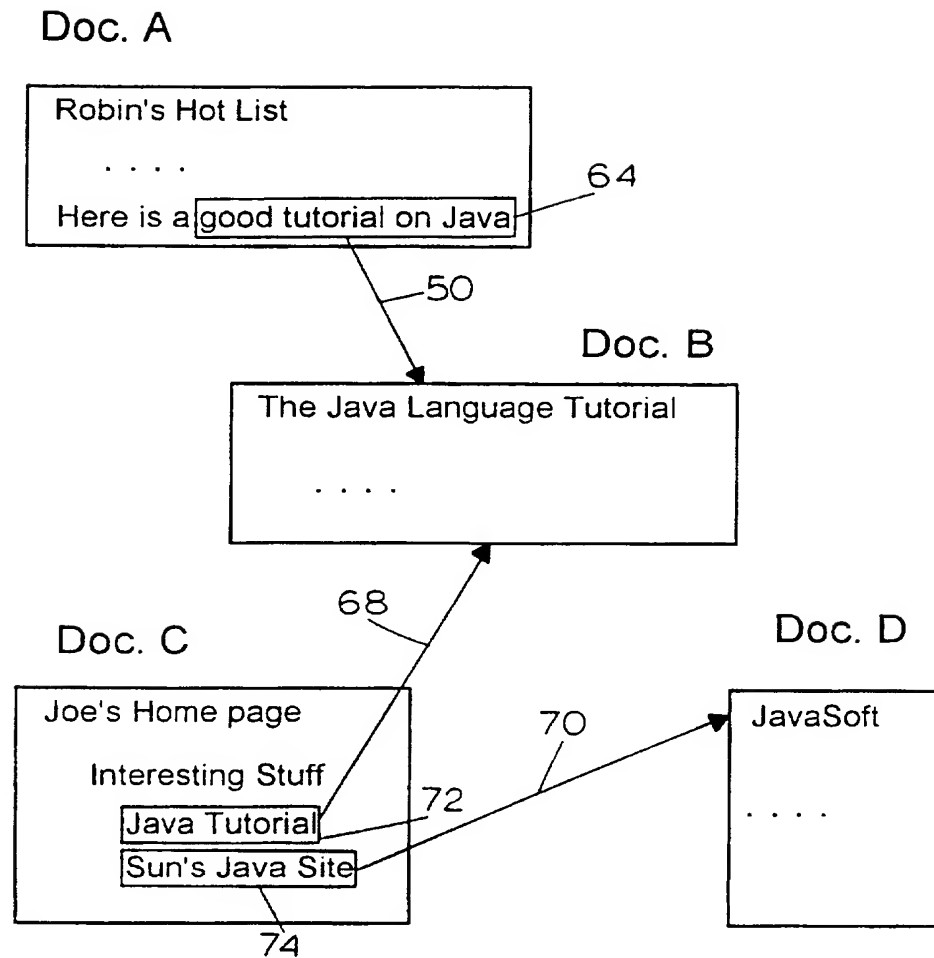
**FIG. 4**

FIG. 5

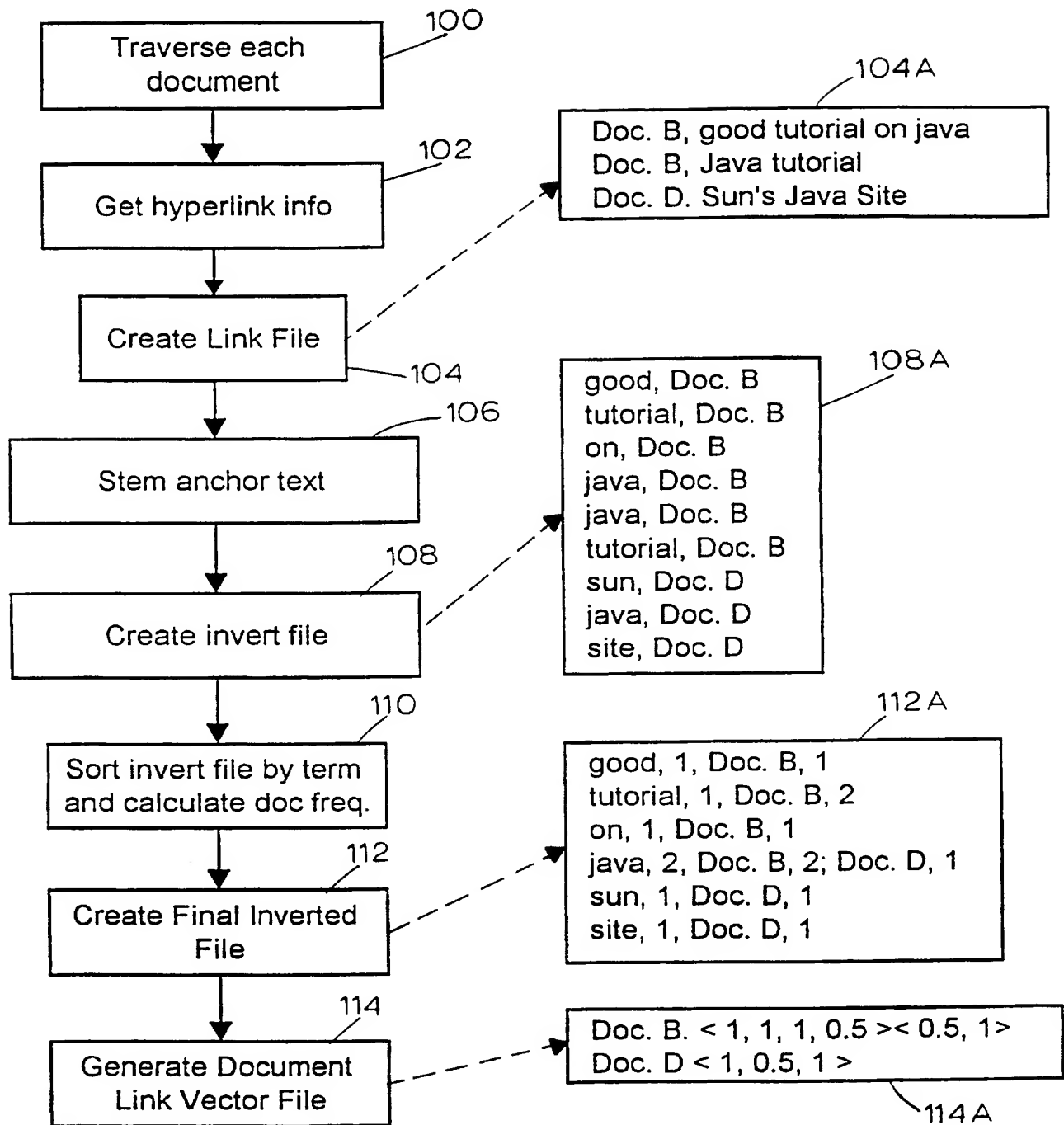
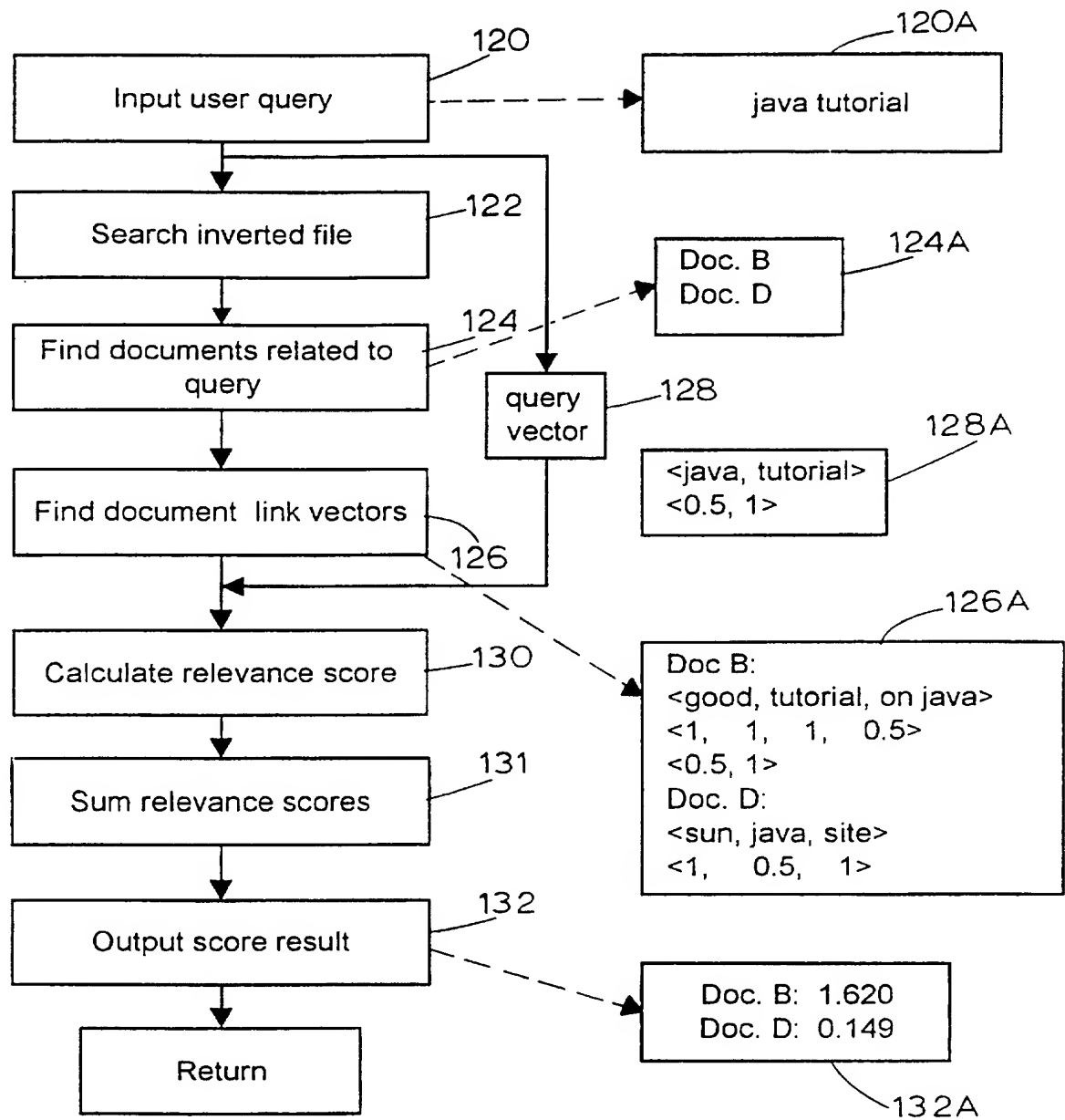


FIG. 6



INTERNATIONAL SEARCH REPORT

International Application No
PCT/US 97/10191

A. CLASSIFICATION OF SUBJECT MATTER
IPC 6 G06F17/30

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)
IPC 6 G06F

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	<p>FREI H P ET AL: "The use of semantic links in hypertext information retrieval" INFORMATION PROCESSING & MANAGEMENT (INCORPORATING INFORMATION TECHNOLOGY), vol. 31, no. 1, January 1995, page 1-13 XP004040956 see abstract see page 2, line 3, paragraph 2.1 - page 3, paragraph 2.3 see page 6, paragraph 4.1 - page 8, paragraph 4.2</p> <p style="text-align: center;">---</p> <p style="text-align: center;">-/-</p>	1-25

☒ Further documents are listed in the continuation of box C.

☐ Patent family members are listed in annex.

* Special categories of cited documents :

- "A" document defining the general state of the art which is not considered to be of particular relevance
- "E" earlier document but published on or after the international filing date
- "L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- "O" document referring to an oral disclosure, use, exhibition or other means
- "P" document published prior to the international filing date but later than the priority date claimed

- "T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
- "X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
- "Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.
- "&" document member of the same patent family

Date of the actual completion of the international search

10 October 1997

Date of mailing of the international search report

20. 10. 97

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Authorized officer

Fournier, C

INTERNATIONAL SEARCH REPORT

International Application No

PCT/US 97/10191

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT

Category °	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	DUNLOP M D ET AL: "Hypermedia and free text retrieval" INFORMATION PROCESSING & MANAGEMENT, 1993, UK, vol. 29, no. 3, ISSN 0306-4573, pages 287-298, XP002043306 see page 289, line 20 - page 290, line 23 ---	1-25
A	BICHTELER J ET AL: "The combined use of bibliographic coupling and cocitation for document retrieval" JOURNAL OF THE AMERICAN SOCIETY FOR INFORMATION SCIENCE, JULY 1980, USA, vol. 31, no. 4, ISSN 0002-8231, pages 278-282, XP002043307 see the whole document -----	1-25